

## Amendments to the Claims

- 1 1. (currently amended) A method for detecting highlights from videos,  
2 comprising:  
3       extracting audio features from the video;  
4       classifying the audio features as labels;  
5       extracting visual features from the video;  
6       classifying the visual features as labels; and  
7       fusing, probabilistically, the audio labels and visual labels into a  
8 single discrete-observation coupled hidden Markov model to detect  
9 highlights in the video.
- 1 2. (currently amended) The method of claim 1, in which the video is  
2 compressed, and the single discrete-observation coupled hidden Markov  
3 model includes the audio features, the visual features, audio states of the  
4 audio features and visual states of the visual features.
- 1 3. (original) The method of claim 1, in which silent features are classified  
2 according to audio energy and zero cross rate.
- 1 4. (original) The method of claim 1, in which the audio features are Mel-  
2 scale frequency cepstrum coefficients.
- 1 5. (original) The method of claim 1, in which the audio features are MPEG-7  
2 descriptors.

- 1 6. (original) The method of claim 1, in which the audio features are  
2 classified using Gaussian mixture models.
- 1 7. (original) The method of claim 1, in which the audio labels are selected  
2 from the group consisting of applause, cheering, ball hit, music, male  
3 speech, female speech, and speech with music.
- 1 8. (original) The method of claim 1, in which the visual features are based  
2 on motion activity descriptors.
- 1 9. (original) The method of claim 1, in which the visual features include  
2 dominant color and motion vectors.
- 1 10. (original) The method of claim 1, in which a variance of the motion  
2 activity is quantized to obtain the visual labels.
- 1 11. (original) The method of claim 1, in which the motion activity is  
2 averaged to obtain the visual labels.
- 1 12. (original) The method of claim 1, in which the visual labels are selected  
2 from the group consisting of close shot, replay, and zoom.
13. (canceled)

1 14. (currently amended) The method of ~~claim 13~~ claim 1, in which the  
2 discrete-observation coupled hidden Markov model includes audio hidden  
3 Markov models and visual hidden Markov models.

1 15. (original) The method of claim 14, in which the discrete-observation  
2 coupled hidden Markov model is generated from a Cartesian product of  
3 states of the audio hidden Markov models and the visual hidden Markov  
4 models, and a Cartesian product of observations of the audio hidden Markov  
5 models and the visual hidden Markov models.

1 16. (currently amended) The method of ~~claim 13~~ claim 1, further  
2 comprising:  
3 training the discrete-observation coupled hidden Markov model with  
4 hand labeled videos.

1 17. (original) The method of claim 1, in which the video is a sport video.

1 18. (original) The method of claim 1, further comprising:  
2 determining likelihoods for the highlights; and  
3 thresholding the highlights.

1 19. (currently amended) The method of claim 2, in which transitions  
2 between the audio states and the visual states of the single discrete-  
3 observation coupled hidden Markov model are represented by transition  
4 matrices ~~portion of the video is compressed~~.

1 20. (currently amended) The method of claim 1, in which the probabilistic  
2 fusion is according to a function fusion function  $F(f_A, f_B)$ , where  $f_A$  are the  
3 audio features and  $f_B$  are the visual features ~~visual portion of the video is~~  
4 ~~compressed.~~

1 21. (new) The method of claim 19, in which the transition matrices have a  
2 form:

$$a_{(i,j),k}^1 = Pr(S_{t+1}^1 = k | S_t^1 = i, S_t^2 = j)$$

$$1 \leq i, k \leq M; 1 \leq j \leq N$$

3 \_\_\_\_\_  
4 where  $S^1$  represents the audio states and  $S^2$  the visual states, and  $Pr$  is a  
5 probability of transitioning to a next state  $k$  at the next time instant  $t$  given  
6 two current hidden states  $i$  and  $j$ , respectively, and  $M$  is a number of audio  
7 states and  $N$  a number of visual states.